

Mirror Descent with Relative Smoothness in Measure Spaces, with application to Sinkhorn and Expectation-Maximization (EM)

Pierre-Cyril Aubin-Frankowski¹, Anna Korba², Flavien Léger³

¹ DI, Ecole normale supérieure, Université PSL, CNRS, INRIA SIERRA, Paris, France

²CREST, ENSAE, IP Paris

³ INRIA MOKAPLAN Paris

NeurIPS 2022

Quick Summary

- Rigorous proof of convergence of Mirror Descent (MD) under relative smoothness and convexity, in the infinite-dimensional setting of optimization over measure spaces
- New and simple way to derive rates of convergence for Sinkhorn's algorithm as an MD over transport plans
- New expression of Expectation-Maximization (EM) as MD, convergence rates when restricted to the latent distribution, coincides with Lucy-Richardson's algorithm in signal processing

Optimisation over the space of measures

Let $\mathcal{X} \subset \mathbb{R}^d$, $\mathcal{M}(\mathcal{X})$ the space of Radon measures on \mathcal{X} , convex functionals $\mathcal{F}, \phi : \mathcal{M}(\mathcal{X}) \rightarrow \mathbb{R} \cup \{+\infty\}$, convex $\mathcal{C} \subset \mathcal{M}(\mathcal{X})$, consider mirror descent:

$$\min_{\mu \in \mathcal{C}} \mathcal{F}(\mu)$$

$$\mu_{n+1} = \operatorname{argmin}_{\nu \in \mathcal{C}} \{d^+ \mathcal{F}(\mu_n)(\nu - \mu_n) + LD_\phi(\nu | \mu_n)\} \quad (1)$$

Under which assumptions does it converge and at which rate?

Examples of optimization of measures

The “Kullback-Leibler divergence” or relative entropy is

$$\text{KL}(\mu|\bar{\mu}) = \begin{cases} \int_{\mathbb{R}^d} \log\left(\frac{\mu}{\bar{\mu}}(x)\right) d\mu(x) & \text{if } \mu \ll \bar{\mu} \\ +\infty & \text{else.} \end{cases}$$

- Entropic optimal transport $\min_{\pi \in \Pi(\mu, \nu)} \text{KL}(\pi|R)$ for $R \propto \exp(-c(x,y)/\epsilon)\mu \otimes \nu$
- Expectation-Maximization $\min_{q \in \mathcal{Q}} \text{KL}(\bar{\nu}|p_y p_q)$ with the observations $\bar{\nu}$
- Bayesian inference $\min_{\mu \in \mathcal{P}(\mathcal{X})} \text{KL}(\mu|\bar{\mu})$ with the posterior $\bar{\mu} \propto \exp(-V)$
- Optimization of 1-hidden layer neural network $\min_{\mu \in \mathcal{C}} \text{MMD}^2(\mu|\bar{\mu})$

Definitions of derivatives

$$\mu_{n+1} = \operatorname{argmin}_{\nu \in \mathcal{C}} \{d^+ \mathcal{F}(\mu_n)(\nu - \mu_n) + LD_\phi(\nu | \mu_n)\}$$

The KL does not have a Gâteaux derivative! Need for weaker notions:

(*directional derivative*) $d^+ \mathcal{F}(\nu)(\mu) = \lim_{h \rightarrow 0^+} \frac{\mathcal{F}(\nu + h\mu) - \mathcal{F}(\nu)}{h},$ (2)

(*first variation*) $\langle \nabla_{\mathcal{C}} \mathcal{F}(\mu), \xi \rangle = d^+ \mathcal{F}(\mu)(\xi) \quad \xi + \mu \in \operatorname{dom}(\mathcal{F}) \cap \mathcal{C},$ (3)

(*Bregman divergence*) $D_\phi(\nu | \mu) = \phi(\nu) - \phi(\mu) - d^+ \phi(\mu)(\nu - \mu).$ (4)

Convergence result for mirror descent under relative smoothness

\mathcal{F} is L -smooth relative to ϕ over C for $L \geq 0$ if, for any $\mu, \nu \in C \cap \text{dom}(\mathcal{F}) \cap \text{dom}(\phi)$,

$$D_{\mathcal{F}}(\nu|\mu) = \mathcal{F}(\nu) - \mathcal{F}(\mu) - d^+\mathcal{F}(\mu)(\nu - \mu) \leq LD_{\phi}(\nu|\mu).$$

Conversely, \mathcal{F} is l -strongly convex relative to ϕ , for $l \geq 0$, if we have

$$D_{\mathcal{F}}(\nu|\mu) \geq lD_{\phi}(\nu|\mu).$$

Theorem 1

Assume that \mathcal{F} is l -strongly convex and L -smooth relative to ϕ , with $l, L \geq 0$. Consider the mirror descent scheme (1), and assume that for each $n \geq 0$, $\nabla_C \phi(\mu_n)$ exists. Then for all $n \geq 0$ and all $\nu \in C \cap \text{dom}(\mathcal{F}) \cap \text{dom}(\phi)$:

$$\mathcal{F}(\mu_n) - \mathcal{F}(\nu) \leq \frac{lD_{\phi}(\nu|\mu_0)}{\left(1 + \frac{l}{L-1}\right)^n - 1} \leq \frac{L}{n}D_{\phi}(\nu|\mu_0)$$

Entropic optimal transport and Sinkhorn

$$\text{Entropic optimal transport } \min_{\pi \in \Pi(\bar{\mu}, \bar{\nu})} \text{KL}(\pi | e^{-c/\epsilon} \bar{\mu} \otimes \bar{\nu})$$

The Sinkhorn algorithm in its primal formulation does alternative (entropic) projections on $\Pi(\bar{\mu}, *)$ and $\Pi(*, \bar{\nu})$, i.e. initializing with $\pi_0 \in \Pi_c$, iterate

$$\pi_{n+\frac{1}{2}} = \operatorname{argmin}_{\pi \in \Pi(\bar{\mu}, *)} \text{KL}(\pi | \pi_n), \quad (5)$$

$$\pi_{n+1} = \operatorname{argmin}_{\pi \in \Pi(*, \bar{\nu})} \text{KL}(\pi | \pi_{n+\frac{1}{2}}). \quad (6)$$

For $c \in L^\infty$, define $C = \Pi(*, \bar{\nu})$ and the objective function $F_S(\pi) = \text{KL}(\rho_{\mathcal{X}} \pi | \bar{\mu})$.

The Sinkhorn iterations can be written as a mirror descent with objective F_S and Bregman divergence KL over the constraint $C = \Pi(*, \bar{\nu})$, with $\nabla F_S(\pi_n) = \ln(d\mu_n/d\bar{\mu}) \in L^\infty(\mathcal{X} \times \mathcal{Y})$, $\mu_n = \rho_{\mathcal{X}} \pi_n$

$$\pi_{n+1} = \operatorname{argmin}_{\pi \in C} \langle \nabla F_S(\pi_n), \pi - \pi_n \rangle + \text{KL}(\pi | \pi_n) \quad (7)$$

Entropic optimal transport and Sinkhorn (cont.)

The functional $F_S(\pi) = \text{KL}(\rho_{\mathcal{X}}\pi|\bar{\mu})$ is convex and is 1-relatively smooth w.r.t. KL over $\mathcal{P}(\mathcal{X} \times \mathcal{Y})$.

$D_C := \frac{1}{2} \sup_{x,y,x',y'} [c(x,y) + c(x',y') - c(x,y') - c(x',y)]$. For $\tilde{\pi}, \pi \in \Pi_C \cap \mathcal{C}$, we have that

$$\text{KL}(\tilde{\pi}|\pi) \leq (1 + 4e^{3D_C/\epsilon}) \text{KL}(\rho_{\mathcal{X}}\tilde{\pi}|\rho_{\mathcal{X}}\pi),$$

i.e. F_S is $(1 + 4e^{3D_C/\epsilon})^{-1}$ -relatively strongly convex w.r.t. KL over $\Pi_C \cap \mathcal{C}$ (cyclically invariant).

For all $n \geq 0$, the Sinkhorn algorithm is a mirror descent and verifies, for π_* the optimum of EOT and μ_* its first marginal,

$$\text{KL}(\mu_n|\mu_*) \leq \frac{\text{KL}(\pi_*|\pi_0)}{(1 + 4e^{\frac{3D_C}{\epsilon}}) \left(\left(1 + 4e^{-\frac{3D_C}{\epsilon}}\right)^n - 1 \right)} \leq \frac{\text{KL}(\pi_*|\pi_0)}{n}.$$

EM and latent EM

We posit a joint distribution $p_q(dx, dy)$ parametrized by an element q of some given set \mathcal{Q} . For $p_y p_q(dy) = \int_{\mathcal{X}} p_q(dx, dy)$, the goal is to infer q by solving

$$\min_{q \in \mathcal{Q}} \text{KL}(\bar{\nu} | p_y p_q), \quad (8)$$

EM then proceeds by alternate minimizations of $\text{KL}(\pi, p_q)$:

$$q_n = \operatorname{argmin}_{q \in \mathcal{Q}} \text{KL}(\pi_n | p_q), \quad (9)$$

$$\pi_{n+1} = \operatorname{argmin}_{\pi \in \Pi(*, \bar{\nu})} \text{KL}(\pi | p_{q_n}). \quad (10)$$

Define the constraint set $\mathcal{C} = \Pi(*, \bar{\nu})$ and $F_{\text{EM}}(\pi) = \inf_{q \in \mathcal{Q}} \text{KL}(\pi | p_q)$.

EM is a mirror descent, with $\nabla F_{\text{EM}}(\pi_n) = \ln(d\pi_n/dp_{q_n})$,

$$\pi_{n+1} = \operatorname{argmin}_{\pi \in \mathcal{C}} \langle \nabla F_{\text{EM}}(\pi_n), \pi - \pi_n \rangle + \text{KL}(\pi | \pi_n) \quad (11)$$

EM and latent EM (cont.)

$F_{\text{EM}} = \inf_{q \in \mathcal{Q}} \text{KL}(\pi | p_q)$ is in general non-convex.

However, writing $p_q(dx, dy) = \mu(dx)K(x, dy)$ and optimizing only over its first marginal, i.e. $q = \mu$, makes F_{EM} convex.

Define $F_{\text{LEM}}(\pi) := \text{KL}(\pi | p_{\mathcal{X}}\pi \otimes K) = \inf_{\mu \in \mathcal{P}(\mathcal{X})} \text{KL}(\pi | \mu \otimes K)$

Latent EM can be written as mirror descent with objective F_{LEM} , Bregman potential ϕ_e and the constraints $\mathcal{C} = \Pi(*, \bar{\nu})$,

$$\pi_{n+1} = \underset{\pi \in \mathcal{C}}{\operatorname{argmin}} \langle \nabla F_{\text{LEM}}(\pi_n), \pi - \pi_n \rangle + \text{KL}(\pi | \pi_n) \quad (12)$$

Set $\mu_* \in \underset{\mu \in \mathcal{P}(\mathcal{X})}{\operatorname{argmin}} \text{KL}(\bar{\nu} | T_K(\mu))$ where $T_K : \mu \in \mathcal{P}(\mathcal{X}) \mapsto \int_{\mathcal{X}} \mu(dx)K(x, \cdot) \in \mathcal{M}(\mathcal{Y})$. The functional F_{LEM} is convex and 1-smooth relative to ϕ_e . For $\pi_0 \in \Pi(*, \bar{\nu})$,

$$\text{KL}(\bar{\nu} | T_K \mu_n) \leq \text{KL}(\bar{\nu} | T_K \mu_*) + \frac{\text{KL}(\mu_* | \mu_0) + \text{KL}(\bar{\nu} | T_K \mu_*) - \text{KL}(\bar{\nu} | T_K \mu_0)}{n}. \quad (13)$$